



# Efficient computation of polynomial explanations of Why-Not questions

Nicole Bidoit, Melanie Herschel, Katerina Tzompanaki

## ► To cite this version:

Nicole Bidoit, Melanie Herschel, Katerina Tzompanaki. Efficient computation of polynomial explanations of Why-Not questions. 31ème Conférence sur la Gestion de Données - Principes, Technologies et Applications - BDA 2015, Sep 2015, Île de Porquerolles, France. hal-01182104

**HAL Id: hal-01182104**

**<https://hal.science/hal-01182104>**

Submitted on 30 Jul 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Efficient computation of polynomial explanations of Why-Not questions

Nicole Bidoit  
Université Paris Sud / Inria  
91405 Orsay Cedex, France  
nicole.bidoit@lri.fr

Melanie Herschel  
Universität Stuttgart  
70569 Stuttgart, Germany  
melanie.herschel  
@ipvs.uni-stuttgart.de

Katerina Tzompanaki  
Université Paris Sud / Inria  
91405 Orsay Cedex, France  
katerina.tzompanaki@lri.fr

## ABSTRACT

Answering a Why-Not question consists in explaining why the result of a query does *not* contain some expected data, called *missing answers*. This paper [6] focusses on processing Why-Not questions following a query-based approach that identifies the culprit *query* components. The first contribution of this paper is a general definition of a Why-Not explanation by means of a *polynomial*. Intuitively, the polynomial provides all possible explanations to explore in order to recover the missing answers. Moreover, this formalism allows us to represent Why-Not explanations to extended relational models having for instance probabilistic or bag semantics. Computing the Why-Not explanation is a complex process and the second contribution of the paper is an algorithm that *efficiently* generates the aforementioned polynomials that answer Why-Not questions. An experimental evaluation demonstrates the practicality of the algorithm both in terms of efficiency and explanation quality, compared to existing algorithms.

Répondre à des questions de type "Pourquoi pas" (Why Not) consiste à expliquer pourquoi certaines données appelées réponses manquantes sont absentes du résultat d'une requête. Cet article traite de questions de type "Pourquoi pas" en suivant une approche "requête", c'est à dire que les explications sont fournies par les combinaisons de conditions de la requête qui sont responsables de la non obtention de certaines réponses. La première contribution est une définition générale de ce qu'est l'explication d'une question "Pourquoi pas" sous la forme d'un polynôme. Intuitivement, ce polynôme fournit toutes les voies à explorer pour récupérer les réponses manquantes. De plus, cette définition permet, avec le même formalisme, de s'intéresser à des extensions du modèle relationnel tel que la sémantique multi-ensembliste ou probabiliste. La deuxième contribution de cet article est liée au calcul des explications d'une question "Pourquoi pas". Un algorithme efficace est présenté, accompagné d'une validation expérimentale et d'une étude comparative.

## 1. INTRODUCTION

The increasing load of data produced nowadays is coupled with an increasing need for complex data transformations that develop-

ers design to process these data in every-day tasks. These transformations, commonly specified declaratively, may result in unexpected outcomes. For instance, given the sample query and data of Fig. 1 on airlines and destination countries, a developer (or traveller) may wonder why Emirates does not appear in the result. Traditionally, she would repeatedly manually analyze the query to identify a possible reason, fix it, and test it to check whether the *missing answer* is now present or if other problems need to be fixed.

Answering such *Why-Not questions*, that is, understanding why some data are *not* part of the result, is very valuable in a series of applications, such as query debugging and refinement, data verification or what-if analysis. To help developers *explain missing answers*, different algorithms have recently been proposed for relational and SQL queries as well as other types of queries like Top-k and reverse skyline.

For relational queries, Why-Not questions can be answered for example based on the data (instance-based explanations), the query (query-based explanations), or both (hybrid explanations). We focus on solutions producing query-based explanations, as these are generally more efficient while providing sufficient information for query analysis and debugging. Essentially, a query-based explanation is a set of conditions of the query that are responsible for pruning data relevant to the missing answers. Existing methods producing query based explanations are not satisfactory, as they return different explanations for the same SQL query, and miss explanations. This is due to the fact that these algorithms are designed over query trees and thus, the explanations depend on the topology of a given tree and indeed to the ordering of the query operators in the query tree.

**EXAMPLE 1.1.** Consider the SQL query and data of Fig. 1 and assume that a developer wants an explanation for the absence of Emirates from the query result. Fig. 2 shows two possible query trees for the query. It also shows the tree operators that Why-Not [9] (◦) and NedExplain [5] (★) return as query-based explanations as well as the tree operators returned as part of hybrid explanations by Conseil [18, 19] (●). Each algorithm returns a different result for each of the two query trees, and in most cases, it is only a partial result as the true explanation of the missing answer is that both the selection is too strict for the tuple (Emirates, 1985, 3) from table Airline and this tuple does not find join partners in table Country.

The above example clearly shows the shortcomings of existing algorithms. Indeed, the developer first has to understand and reason at the level of query trees instead of reasoning at the level of the declarative SQL query she is familiar with. Second, she always has to wonder whether the explanation is complete, or if there are other explanations that she could consider instead. To this problem, we make in this paper the following contributions:

(c) 2015, Copyright is with the authors. Informally presented at the BDA 2015 Conference (September 29-October 2, 2015, Ile de Porquerolles, France).

(c) 2015, Droits restant aux auteurs. Présenté à la conférence BDA 2015 (29 Septembre-02 Octobre 2015, Ile de Porquerolles, France).

SELECT airline, country FROM Airline A, Country C WHERE ccode = code AND year < 1985	Airline			Country	
	airline	year	ccode	code	country
	KLM	1919	1	1	Australia
	Qatar	1993	1	2	France
	Aegean	1987	2		
	Emirates	1985	3		

Figure 1: Example query and data

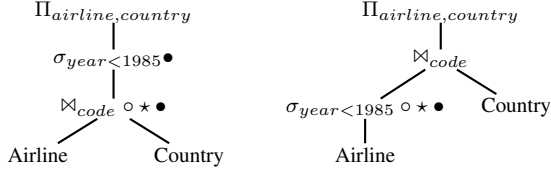


Figure 2: Reordered query trees for the query of Fig. 1 and algorithm results (Why-Not  $\circ$ , NedExplain  $\star$ , Conseil  $\bullet$ )

### Extended formalization of Why-Not explanation polynomial.

In preliminary work [3, 4], we introduced polynomials as Why-Not explanations in the context of the relational model under set semantics. A polynomial provide a complete explanation and is independent of a specific query tree representation, solving the problems illustrated by the Ex. 1.1.

This paper significantly extends on the preliminary notion of a Why-Not explanation: the overall framework has been considerably simplified while the notion of Why-Not explanation is extended to be used in the context of relational model under set, bag and probabilistic semantics. This confirms the robustness of the chosen polynomial representation, making it a good fit for a unified framework for representing Why-Not explanations.

### Efficient Ted++ algorithm.

In our preliminary work [3, 4], we presented a naive algorithm computing Why-Not explanations. We show that its runtime complexity is impractical and propose a totally new algorithm, *Ted++*. *Ted++* is capable of efficiently computing the Why-Not explanation polynomial, based on techniques like smart data partitioning (allowing for a distributed computation) and advantageous translation of expensive database evaluations by mathematical calculations.

### Experimental validation.

We experimentally validate both the efficiency and the effectiveness of the solutions proposed in this paper. These experiments include a comparative evaluation to existing algorithms computing query-based explanations for SQL queries (or sub-languages thereof) as well as a thorough study of *Ted++* performance w.r.t. different parameters.

The remainder of this paper is structured as follows. Sec. 2 covers related work. Sec. 3 defines in detail our problem setting and the Why-Not explanation polynomials. Next, we discuss in detail the *Ted++* algorithm in Sec 4. Finally, we present our experimental setup and evaluation in Sec. 5 before we conclude in Sec. 6.

## 2. RELATED WORK

Recently, we observe the trend that growing volumes of data are processed by programs developed not only by expert developers but also by less knowledgeable users (creation of mashups, use of web services, etc.). These trends have led to the necessity of providing algorithms and tools to better understand and verify the behavior and semantics of developed data transformations, and various so-

Table 1: Algorithms for answering Why-Not questions

Algorithm	Why-Not question	Explanation format	Query
<b>Query-based explanations</b>			
Why-Not [9]	simple	query operators	SPJU
NedExplain [5]	simple	query operators	SPJUA
Ted [3]/Ted++	complex	polynomial	conj. queries with inequalities
<b>Hybrid explanations</b>			
Conseil [18, 19]	simple	source table edits + query operators	SPJAN
<b>Instance-based explanations</b>			
MA [21]	simple	source table edits	SPJ
Artemis [20]	complex	source table edits	SPJUA
Meliou <i>et. al.</i> [24]	simple	causes (tuples) and responsibility	conjunctive queries
Calvanese <i>et. al.</i> [7]	simple	additions to ABox	instance & conj. queries over DL-Lite ontology
<b>Ontology-based explanations</b>			
Cate <i>et. al.</i> [8]	simple	tuples concepts	conj. queries with comparisons
<b>Refinement-based explanations</b>			
ConQueR [28]	complex	rewritten query	SPJA
Zhang <i>et. al.</i> [17]	simple	refined query	Top-k query
Islam <i>et. al.</i> [22]	simple	refined query & Why-Not question	Reverse skyline query
WQRTQ [13]	simple	refined query & Why-Not question	Reverse Top-k query
Chen <i>et. al.</i> [10]	simple	refined query	Spatial keyword Top-k query

lutions have been proposed so far, including data lineage [12] and more generally data provenance [11], (sub-query) result inspection and explanation [2, 15, 27], query conditions relaxation [25], transformation specification simplification [23, 26], etc.

The work presented in this paper falls in the category of data provenance research, and specifically explaining missing answers from query results. Due to the lack of space, the subsequent discussion focuses on this sub-problem, thus on algorithms answering Why-Not questions. Tab. 1 summarizes these algorithms, first classifying them according to the type of explanation they generate (instance-based, query-based, hybrid, ontology-based or refinement-based). The table further shows whether an algorithm supports *simple* Why-Not questions, i.e., questions where each condition impacts one relation only, or more *complex* ones. The last two columns summarize the form of a returned explanation and the class of query supported by an algorithm respectively.

### Query-based and hybrid explanations.

Why-Not [9] takes as input a simple Why-Not question and returns so called *picky query operators* as query-based explanation. To determine these, the algorithm first identifies tuples in the source database that satisfy the conditions of the input Why-Not question and that are not part of the lineage [12] of any tuple in the query result. These tuples, named *compatible tuples*, are traced through the query operators of a query tree representation to identify which operators include them in their input but not in their output. In [9] the algorithm is shown to work for queries involving selection, projection, join, and union (SPJU query). NedExplain [5] is very similar to Why-Not in the sense that it supports simple Why-Not questions and returns a set of *picky operators* as query-based Why-Not explanation as well. However, it supports a broader range of queries, i.e., queries involving selection, projection, join, and aggregation (SPJA queries) and unions thereof and the computation of *picky operators* is significantly different. In this work, we support a wider class of Why-Not questions (complex ones) and provide a new formalization of Why-Not explanation as polynomials.

Conseil [18, 19] produces hybrid explanations that include an instance-based and a query-based component. The latter consists in a set of *picky query operators*. However, as Conseil considers both the data to be possibly incomplete and the query to be possi-

bly faulty, the set of picky query operators associated to a hybrid explanation depends on the set of source edits of the same hybrid explanation.

### Instance-based explanations.

Both Missing-Answers (MA) [21] and Artemis [20] compute instance-based explanations in the form of source table edits. Whereas MA returns correct explanations for simple Why-Not questions and SQL queries involving selection, projection, and join (SPJ queries), Artemis supports complex Why-Not questions on a larger fraction of SQL queries (including union or aggregation, denoted SPJUA). Meliou *et al.* [24] study the unification of instance-based explanations of missing answers and of data *present* in a conjunctive query result, leveraging the concepts of causality and responsibility. Finally, Calvanese *et al.* [7] leverage abductive reasoning and theoretically examines the problem of computing instance-based explanations for a class of simple Why-Not questions on data represented by a DL-Lite ontology. As here the explanations are in the form of source edits, we consider these works orthogonal to query-based ones.

### Ontology-based explanations.

Cate *et al.* [8] introduce ontology-based explanations for conjunctive queries with comparisons ( $=$ ,  $>$ ,  $<$ ), using external or data workload generated ontologies. They also provide an algorithm that computes these ontology-based explanations. The algorithm as well as the explanations are completely independent of the query to be analysed and therefore, we consider this approach orthogonal to our work here.

### Refinement-based explanations.

Given a set of missing answers and a SPJUA query, ConQueR [28] refines the query such that all missing answers become part of the output. Refinements of the query and of the Why-Not question have been proposed in other contexts as well like for Top-k queries (Zhang *et al.* [17]), reverse skyline queries (Islam *et al.* [22]), reverse Top-k queries (WQRTQ [13]), spatial keyword Top-K queries (Chen *et al.* [10]) etc. Although these approaches are generally very interesting, they do not focus on pinpointing to the user the erroneous parts of the query, but on directly refining the query. Indeed, the generated queries may contain changes that are not necessarily tied to an erroneous part of the query. For this reason, they are out of the scope of this paper.

## 3. WHY-NOT EXPLANATION POLYNOMIAL

This section introduces a polynomial formalization of query-based Why-Not explanations. We assume the reader familiar with the relational model [1], and we only briefly revisit some relevant notions in Sec. 3.1 while we formalize Why-Not questions. Then, in Sec. 3.2, we define the explanation of a Why-Not question as a polynomial and in Sec. 3.3 we provide a unified general framework for Why-Not explanations in the context of probabilistic, set, and bag semantics databases.

### 3.1 Preliminaries

For the moment, we limit our discussion to relational databases under set semantics. A database schema  $\mathcal{S}$  is a set of relation schemas. A relation schema  $R$  is a set of attributes. We assume each attribute of  $R$  qualified, i.e., of the form  $R.A$  and, for the sake of simplicity we assume a unique domain  $D$ .  $\mathcal{I}$  denotes a database instance over  $\mathcal{S}$  and  $\mathcal{I}_R$  denotes the instance of a relation  $R \in \mathcal{S}$ . We assume that each database relation  $R$  has a special attribute  $R.Id$ ,

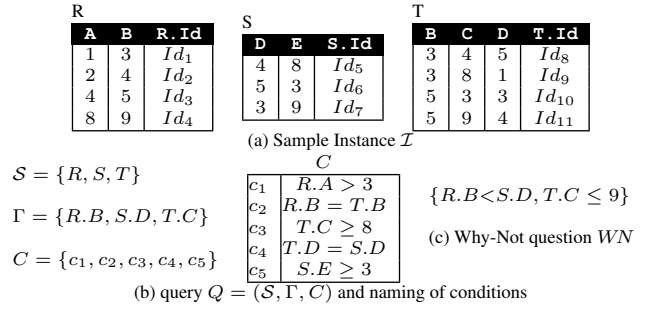


Figure 3: Running example

which is used as identifier for the tuples in  $\mathcal{I}_R$ . For any object  $O$  (a relational schema, database schema, condition etc),  $\mathcal{A}(O)$  denotes the set of attributes occurring in  $O$ . Finally, a condition  $c$  over  $\mathcal{S}$  is defined as an expression of the form  $R.A \theta a$  where  $a \in D$  or of the form  $R.A \theta S.B$ , where  $R.A, S.B \in \mathcal{A}(\mathcal{S})$ , and  $\theta \in \{=, \neq, <, \leq\}$ . A condition over two relations is *complex*, otherwise it is *simple*. In this article, we consider conjunctive queries with inequalities. Note that, in our approach, the database schema  $\mathcal{S}$  denotes the query input schema. In an SQL-like approach, each time we need an instance of a relation, we refer to it by a different name. In this way, we are able to correctly define Why-Not questions in the case of self-join.

**DEFINITION 3.1 (QUERY Q).** A query  $Q$  is specified by the triple  $(\mathcal{S}, \Gamma, C)$ , where  $\mathcal{S}$  is a database schema,  $\Gamma \subseteq \mathcal{A}(\mathcal{S})$  is the projection attribute set, and  $C$  is a set of conditions over  $\mathcal{A}(\mathcal{S})$ . The semantics of  $Q$  are given by the relational algebra expression  $\pi_{\Gamma} [\sigma_{\bigwedge_{c \in C} c} [\times_{R \in \mathcal{S}} [R]]]$ .

The result of  $Q$  over a database instance  $\mathcal{I}$  of  $\mathcal{S}$  is denoted by  $Q[\mathcal{I}]$ . Note here that we are not concerned about the evaluation/optimization of  $Q$ .

**EXAMPLE 3.1.** Fig. 3 describes our running example. Fig. 3(a) displays an instance  $\mathcal{I}$  over  $\mathcal{S} = \{R, S, T\}$ . Fig. 3(b) displays a query  $Q$  over  $\mathcal{S}$ , whose conditions have been named for convenience.  $R.B = T.B$  and  $T.D = S.D$  are complex whereas the others are simple conditions. Moreover,  $Q[\mathcal{I}] = \{(R.B:5, S.D:4, T.C:9)\}$ .

In our framework, a Why-Not question specifies missing tuples from the result of a query  $Q$  through a conjunctive set of conditions. As a Why-Not question is related to the result of  $Q$ , the conditions of the Why-Not question are restricted to the attributes of the output schema of  $Q$ .

**DEFINITION 3.2 (WHY-NOT QUESTION).** A Why-Not question  $WN$  w.r.t.  $Q$  is defined as a set of conditions over  $\Gamma$ .

The notion of complex and simple conditions is extended to complex and simple Why-Not questions in a straight forward manner.

As we said, a Why-Not question  $WN$  summarizes a set of (missing) tuples that the user expected to find in the query result. To be able to obtain these missing tuples as query results, data from the input relation instances that satisfy  $WN$  need to be combined by the query. The candidate data combinations are what we call *compatible* tuples and these can be computed using  $WN$  as in Def. 3.3.

**DEFINITION 3.3 (COMPATIBLE TUPLES).** Consider the query  $Q_{WN} = (\mathcal{S}, \mathcal{A}(\mathcal{S}), WN)$ , where  $\mathcal{S}$  is the input schema of  $Q$ . The set  $CT$  of compatible tuples is the result of the query  $Q_{WN}$  over  $\mathcal{I}$ .

We further introduce the notion of a *well founded* Why-Not question. Intuitively, a Why-Not question can be answered under a query-based approach, only if some data in  $\mathcal{I}$  match the Why-Not question (otherwise instance-based explanations should be sought for). Moreover, a Why-Not question is meaningful if it tracks data not already returned by the query.

**DEFINITION 3.4 (WELL FOUNDED WHY-NOT QUESTION).** *A Why-Not question  $WN$  is said to be well founded if  $CT \neq \emptyset$  and  $\pi_T[CT] \cap Q[\mathcal{I}] = \emptyset$ .*

**EXAMPLE 3.2.** *Continuing Ex. 3.1, we may wonder why there is not a tuple for which  $R.B < S.D$  and  $T.C \leq 9$ . According to Def. 3.2, this Why-Not question is the conjunction of the conditions  $R.B < S.D \wedge T.C \leq 9$  (Fig. 3(c)). Since  $R.B < S.D$  is a complex condition,  $WN$  is a complex Why-Not question. The compatible tuples set  $CT$  is the result of the query  $Q_{WN} = \sigma_{R.B < S.D \wedge T.C \leq 9}[R \times S \times T]$ , which contains 12 tuples. For example, one compatible tuple is  $\tau_1 = (R.Id:1, R.A:1, R.B:3, S.Id:5, S.D:4, S.E:8, T.Id:8, T.B:3, T.C:4, T.D:5)$ .*

Each tuple in  $CT$  could have led to a missing tuple, if it was not eliminated by some of the query's conditions. Thus, explaining  $WN$  amounts to identifying these blocking query conditions.

### 3.2 The Why-Not Explanation Polynomial

To build the query-based explanation of  $WN$ , we start by specifying what explains that a compatible tuple  $\tau$  did not lead to an answer. Intuitively, the explanation consists of the query conditions pruning  $\tau$ .

**DEFINITION 3.5 (EXPLANATION FOR  $\tau$ ).** *Let  $\tau \in CT$  be a compatible tuple w.r.t.  $WN$ , given  $Q$ . Then, the explanation for  $\tau$  is the set of conditions  $\mathcal{E}_\tau = \{c | c \in C \text{ and } \tau \not\models c\}$ .*

**EXAMPLE 3.3.** *Reconsider the compatible tuple  $\tau_1$  in Ex. 3.2. The conditions of  $Q$  (see Ex. 3.1), not satisfied by  $\tau_1$  are:  $c_1$ ,  $c_3$ , and  $c_4$ . So, the explanation for  $\tau_1$  is  $\mathcal{E}_{\tau_1} = \{c_1, c_3, c_4\}$ .*

Having defined the explanation wrt one compatible tuple, the explanation for  $WN$  is obtained by simply summing up the explanations for all the compatible tuples in  $CT$ . This leads to an expression of the form  $\sum_{\tau \in CT} \prod_{c \in \mathcal{E}_\tau} c$ . We justify modelling the explanation of  $\tau$  with a product (meaning conjunction) of conditions by the fact that in order for  $\tau$  to 'survive' the query conditions and give rise to a missing tuple, every single condition in the explanation must be 'repaired'. The sum (meaning disjunction) of the products for each  $\tau \in CT$  implies that if any explanation is 'correctly repaired', the associated  $\tau$  will produce a missing tuple.

Of course, several compatible tuples can share the same explanation. Thus, the final Why-Not explanation is a polynomial having as variables the query conditions and as integer coefficients the number of compatible tuples sharing an explanation.

**DEFINITION 3.6 (WHY-NOT EXPLANATION).** *With the same assumption as before, the Why-Not explanation for  $WN$  is defined as the polynomial*

$$PEX = \sum_{\mathcal{E} \in E} \text{coef}_{\mathcal{E}} \prod_{c \in \mathcal{E}} c$$

where  $E = 2^C$ ,  $\text{coef}_{\mathcal{E}} \in \{0, \dots, |CT|\}$  is the number of tuples in  $CT$  sharing  $\mathcal{E}$  as an explanation, and  $\sum_{\mathcal{E} \in E} \text{coef}_{\mathcal{E}} = |CT|$ .

Intuitively,  $E$  contains all possible explanations, i.e., condition combinations and each of these explanations prunes from zero to at most  $|CT|$  compatible tuples. Each compatible tuple is pruned by exactly one condition combination, which is why the sum of all coefficients is equal to the number of compatible tuples.

We mentioned before that each term of the polynomial provides an alternative explanation to be explored by the user who wishes to recover some missing tuples. Additionally, the polynomial as in Def. 3.6 offers through its coefficients some useful hints to users interested in the *number* of recoverable tuples. More precisely, by isolating an explanation  $\mathcal{E}$  to repair, we can obtain an upper bound for the number of compatible tuples that can be recovered. The upper bound is calculated as the sum of the coefficients of all the explanations that are sub-sets of (the set of conditions of)  $\mathcal{E}$ , because when  $\mathcal{E}$  is changed it is likely that some sub-combinations are also repaired.

**EXAMPLE 3.4.** *In Ex. 3.3 we found the explanation  $\{c_1, c_3, c_4\}$ , leading to the polynomial term  $c_1 * c_3 * c_4$ . Taking into consideration all the 12 compatible tuples of our example, we obtain the following PEX polynomial:  $2 * c_1 * c_4 + 2 * c_1 * c_3 * c_4 + 4 * c_1 * c_2 * c_4 + 2 * c_1 * c_2 * c_3 + 2 * c_1 * c_2 * c_3 * c_4$ . In the polynomial, each addend, composed by a coefficient and an explanation, captures a way to obtain missing tuples. For instance, the explanation  $c_1 * c_2 * c_4$  indicates that we may recover some missing answers if  $c_1$  and  $c_2$  and  $c_4$  are changed. Then, the sum of its coefficient 4 and the coefficient 2 of the explanation  $c_1 * c_4$  ( $\{c_1, c_4\} \subseteq \{c_1, c_2, c_4\}$ ) indicates that we can recover from 0 to 6 tuples.*

As the presentation of the polynomial per se may sometimes be cumbersome, and thus not easy for a user to manipulate, some post-processing steps could be applied. For example, depending on the application or needs, only a subset of the explanations could be returned like minimum explanations (i.e., for which no sub-explanations exist), or those explanations supposed to recover a specific number of tuples, or having specific condition types etc.

### 3.3 Extension: Bag & Probabilistic Semantics

So far, we have considered databases under *set* semantics only. In this section, we discuss how the definition of Why-Not explanation (Def. 3.6) extends to settings with conjunctive queries over bag semantics and probabilistic databases.

$K$ -relations, as introduced in [14], capture in a unified manner relations under probabilistic, bag or set semantics. Briefly, a  $K$ -relation maps tuples to elements of a set  $K$ , that is  $K$ -relation tuples are annotated with elements in  $K$ . In our case, we consider that  $K$  is a set of tuple identifiers, similar to our special attribute  $R.Id$  in Sec. 3.1.

In what follows, we use the notion of *how-provenance* of tuples in the result of a query  $Q$ . The how-provenance of  $t \in Q(\mathcal{I})$  is modelled as the polynomial obtained by the positive algebra on  $K$ -relations, proposed in [14]. Briefly, each  $t$  is annotated with a polynomial whose variables are tuple identifiers and coefficients are natural numbers. Following [14]'s algebra, if  $t$  results from a selection operator on a tuple  $t_1$  annotated with  $Id_1$ , then  $t$  is also annotated with  $Id_1$ . If  $t$  is the result of the join of  $t_1$  and  $t_2$ , then  $t$  is annotated with  $Id_1 Id_2$ .

We compute the *generalized* Why-Not explanation polynomial as follows. Firstly, we compute the how-provenance for compatible tuples in  $CT$  by evaluation of the query  $Q_{WN}$  (Def. 3.3) wrt the algebra in [14]. Recall that  $Q_{WN}$  contains only selection and join operators. Thus, we assume that each compatible tuple  $\tau$  in  $CT$  is annotated with its how-provenance polynomial, denoted by  $\eta_\tau$ .

In a second step, we associate the expressions of *how* and *why*-not provenance. In order to do this, for each compatible tuple  $\tau$  in  $CT$ , we combine its how-provenance polynomial  $\eta_\tau$  with its explanation  $\mathcal{E}_\tau$  (Def. 3.5). So, each  $\tau$  is annotated with the expression  $\eta_\tau \mathcal{E}_\tau$ .

Finally, as before, we sum the combined expressions for all compatible tuples. The result is the generalized Why-Not explanation  $PEX_{gen} = \sum_{\tau \in CT} \eta_\tau \mathcal{E}_\tau$ .

We now briefly comment on how  $PEX_{gen}$  is instantiated to deal either with the set, bag or probabilistic semantics. Indeed, the ‘specialization’ of  $PEX_{gen}$  relies on the interpretation of the elements in  $K$ , that is on a function  $Eval$  from  $K$  to some set  $L$ . For the set semantics, each tuple in a relation occurs only once. This results in choosing  $L$  to be the singleton  $\{1\}$  and mapping each tuple identifier to 1. It is then quite obvious to note, for the set semantics, that  $PEX_{gen} = PEX$  (Def. 3.6). In the same spirit, for bag semantics,  $L$  is chosen as the set of natural numbers  $\mathbb{N}$  and each tuple identifier is mapped to its number of occurrences. Finally, for probabilistic databases,  $L$  is chosen as the interval  $[0, 1]$  and each tuple identifier is mapped to its occurrence probability.

Thus, the generalized definition of Why-Not explanation is parameterized by the mapping  $Eval$  of the annotations (elements in  $K$ ) in the set  $L$ .

**DEFINITION 3.7. (Generalized Why-Not explanation polynomial)** Given a query  $Q$  over a database schema  $S$  of  $K$ -relations, the generalized Why-Not explanation polynomial for  $WN$  is

$$PEX_{gen} = \sum_{\mathcal{E} \in E} \sum_{\tau \in CT \text{ s.t. } \mathcal{E}_\tau = \mathcal{E}} Eval(\eta_\tau) \mathcal{E}$$

where  $E=2^C$ ,  $\eta_\tau$  is the how-provenance of  $\tau$ , and  $Eval:K \rightarrow L$  evaluates the elements of  $K$  to values in  $L$ .

The specializations of  $PEX_{gen}$  share the same explanations (terms of the polynomial), capturing the same ‘erroneous’ parts of the query. However, the coefficients are interpreted wrt the how-provenance.

## 4. TED++ ALGORITHM

In [3], we have introduced *Ted*, a naive algorithm that implements the definitions of [3] for Why-Not explanations in a straightforward manner. Briefly, *Ted* enumerates the set of compatible tuples  $CT$  by executing the query  $Q_{WN}$  (Def. 3.3). Then, it computes the explanation for each compatible tuple in  $CT$ , which leads to the computation of the final Why-Not explanation. However, both of these steps make *Ted* computationally prohibitive. Not only is the computation of  $CT$  time and space consuming as it often requires cross product executions, but also the iteration over this (potentially very large) set is time consuming. *Ted*’s time complexity is  $O(n^{|S|})$ ,  $n = \max(\{|\mathcal{I}_R|\}, R \in S)$ . As experiments in Sec. 5 confirm, this complexity renders *Ted* of no practical interest.

To overcome the poor performance of *Ted*, we propose *Ted++*. The main feature of *Ted++* is to completely avoid enumerating and iterating over the set  $CT$ , thus it significantly reduces both space and time consumption. Instead, *Ted++* opts for (i) iterating over the space of possible explanations, which is expected to be much smaller, (ii) computing *partial* sets of *passing* compatible tuples, and (iii) computing the *number of eliminated* compatible tuples for each explanation. Intuitively, passing tuples w.r.t. an explanation are tuples satisfying the conditions of the explanation. Finally, the polynomial is computed based on mathematical calculations.

Theorem 4.1 states that *Ted++* is sound and complete w.r.t. Def. 3.6.

### Algorithm 1: *Ted++*

---

**Input:**  $Q=(S, \Gamma, C), \mathcal{I}, WN$   
**Output:**  $PEX$

```

1  $E \leftarrow \text{powerset}(C);$ 
2  $\mathcal{P} \leftarrow \text{validPartitioning}(S, WN);$  *(Def. 4.1)*
3 for  $Part$  in  $\mathcal{P}$  do
4    $CT_{|Part} \leftarrow (Part, \mathcal{A}(Part), WN_{|Part})[\mathcal{I}_{|Part}];$ 
5  $\text{coefficientEstimation}(E, Partition);$ 
6  $PEX \leftarrow \text{post-processing}();$ 
7 return  $PEX;$ 
```

---

**THEOREM 4.1.** Given a query  $q$ , a Why-Not question  $WN$  and an input instance  $\mathcal{I}$ , *Ted++* computes exactly  $PEX$ .

Alg. 1 provides an outline of *Ted++*. The input includes the query  $Q=(S, \Gamma, C)$ , the Why-Not question  $WN$  and the input instance  $\mathcal{I}$ . Firstly, in Alg. 1, line 1, all potential explanations (combinations of the conditions in  $C$ ) are enumerated ( $E=2^C$ ). The remaining steps, discussed in the next subsections, aim at computing the coefficient of each explanation. To illustrate the concepts introduced in the detailed discussions, we will rely on our running example, for which Fig. 4 shows all relevant intermediate results. It should be read bottom-up. For convenience, in our examples, we use subscript  $i$  instead of  $c_i$ .

### 4.1 Partial Compatible Tuples Computation

Using the conditions in  $WN$ , *Ted++* partitions the schema  $S$  (Alg. 1 line 2) into components of relations connected by the conditions in  $WN$  (Def. 4.1).

**DEFINITION 4.1. (Valid Partitioning of  $S$ ).** Given  $WN$ , the partitioning of a database schema  $S$  into  $k$  partitions, denoted  $\mathcal{P} = \{Part_1, \dots, Part_k\}$ , is valid if each  $Part_i$ ,  $i \in \{1, \dots, k\}$  is minimal w.r.t. the following property:  
if  $R \in Part_i$  and  $R' \in S$  s.t.  $\exists c \in WN$  with  $\mathcal{A}(c) \cap \mathcal{A}(R') \neq \emptyset$  and  $\mathcal{A}(c) \cap \mathcal{A}(R) = \emptyset$  then  $R' \in Part_i$ .

The partitioning of  $S$  allows for handling compatible tuples more efficiently, by ‘cutting’ them in distinct meaningful ‘chunks’, avoiding combining chunks over distinct partitions through cross product. We refer to the chunks of compatible tuples as *partial* compatible tuples and group them in sets depending on the partition they belong to. The set  $CT_{|Part}$  of partial compatible tuples wrt  $Part \in \mathcal{P}$  is obtained by evaluating the query  $Q_{Part} = (Part, \mathcal{A}(Part), WN_{|Part})$  over  $\mathcal{I}_{|Part}$  (Alg. 1, line 4).  $WN_{|Part}$  and  $\mathcal{I}_{|Part}$  denote the restriction of  $WN$  and  $\mathcal{I}$  over the relations in  $Part$ , respectively.

**EXAMPLE 4.1.** The valid partitioning of  $S$  is  $Part_1 = \{R, S\}$  (because of the condition  $R.B < S.D$ ) and  $Part_2 = \{T\}$ . The sets of partial compatible tuples  $CT_{|Part_1}$  and  $CT_{|Part_2}$  are given in the bottom line of Fig. 4.

It is easy to prove that the valid partitioning of  $S$  is unique and that the set  $CT$  can be computed from the individual  $CT_{|Part_i}$ .

**LEMMA 4.1.** Let  $\mathcal{P}$  be the valid partitioning of  $S$ . Then,  
 $CT = \times_{Part_i \in \mathcal{P}} [CT_{|Part_i}]$ .

Indeed, Lemma. 4.1 makes it clear how  $CT$  is computed from partial compatible tuples. Our algorithm is designed in a way that avoids computing  $CT$  and relies on the computation of  $CT_{|Part_i}$  only.

---

**Algorithm 2:** coefficientEstimation
 

---

**Input:**  $E$  explanations space,  $\mathcal{P}$  valid partitioning of  $S$

```

1 for  $\mathcal{E} \in E$  *access in ascending size order* do
2   Compute  $part_{\mathcal{E}}$ ;
3   if  $|\mathcal{E}| = 1$  then
4     materialize  $V_{\mathcal{E}}$ ;
5      $\beta_{\mathcal{E}} \leftarrow \text{Eq. (B)}$ ;
6   else
7     if  $\alpha_{\text{subcombination of } \mathcal{E}} \neq 0$  then
8        $\{\mathcal{E}_1, \mathcal{E}_2\} \leftarrow \text{subCombinationsOf}(\mathcal{E})$ ;
9        $\Gamma_{12} \leftarrow \Gamma_1 \cap \Gamma_2$ ;  $*\Gamma_i$  is the output schema of  $V_{\mathcal{E}_i}$ *
10      if  $\Gamma_{12} \neq \emptyset$  then
11         $V_{\mathcal{E}} \leftarrow V_{\mathcal{E}_1} \bowtie_{\Gamma_{12}} V_{\mathcal{E}_2}$ ;
12        materialize  $V_{\mathcal{E}}$ ;
13      else
14         $|V_{\mathcal{E}}| \leftarrow |V_{\mathcal{E}_1}| * |V_{\mathcal{E}_2}|$ ;
15      else
16         $|V_{\mathcal{E}}| \leftarrow |V_{\mathcal{E}_1}| * |V_{\mathcal{E}_2}|$ ;
17       $\beta_{\mathcal{E}} \leftarrow \prod_{Part \in part_{\mathcal{E}}} |CT|_{Part}| - |(\bigcup_{i=1}^n V_{c_i})^{ext}|$ ;  $*\text{Eq. (E)}*$ 
18     $\alpha_{\mathcal{E}} \leftarrow \text{Eq. (A)}$ ;
  
```

---

Next, we compute the number of compatible tuples eliminated by each possible explanation, starting with the partial compatible tuple sets previously defined. These numbers approximate the coefficient of the explanations in the polynomial. Since from this point on, we are only handling (partial) compatible tuples, we omit the word ‘compatible’ to lighten the discussion.

## 4.2 Polynomial Coefficient Estimation

Each set  $\mathcal{E}$  in the powerset  $E$  is in fact a potential explanation that is further processed. In order to do that, we associate  $\mathcal{E}$  with (i) the set of partitions  $part_{\mathcal{E}}$  on which  $\mathcal{E}$  is defined, (ii) the view definition  $V_{\mathcal{E}}$  meant to store the passing partial tuples wrt  $\mathcal{E}^1$ , and, (iii) the number  $\alpha_{\mathcal{E}}$  of tuples eliminated by  $\mathcal{E}$ .

Alg. 2 describes how we process  $E$  in ascending order of explanation size, in order to compute  $\alpha_{\mathcal{E}}$ . Each step deals with explanations of size  $s$ , reusing results from previous steps avoiding cross product computations through mathematical calculations.

We first determine the set of partitions for an explanation  $\mathcal{E}$  as  $part_{\mathcal{E}} = \bigcup_{c \in \mathcal{E}} \{Part_c\}$ , where  $Part_c$  contains at least one relation over which  $c$  is specified.

**EXAMPLE 4.2.** Consider  $\mathcal{E}_1 = \{c_1\}$  and  $\mathcal{E}_2 = \{c_2\}$ . From Fig. 3(b) and the partitions in Fig. 4, we can see that  $c_1$  impacts only  $Part_1$ , whereas  $c_2$  spans over  $Part_1$  and  $Part_2$ . Hence,  $part_{\mathcal{E}_1} = \{Part_1\}$  and  $part_{\mathcal{E}_2} = \{Part_1, Part_2\}$ . Then,  $\mathcal{E} = \{c_1, c_2\}$  is impacted by the union of  $part_{\mathcal{E}_1}$  and  $part_{\mathcal{E}_2}$ , thus  $part_{\mathcal{E}} = \{Part_1, Part_2\}$ .

We use Eq. (A) to calculate the number  $\alpha_{\mathcal{E}}$  of eliminated tuples, using the number  $\beta_{\mathcal{E}}$  of eliminated partial tuples and the cardinality of the partitions not in  $part_{\mathcal{E}}$ . Intuitively, this formula extends the partial tuples to “full” tuples over  $CT$ ’s schema.

$$\alpha_{\mathcal{E}} = \beta_{\mathcal{E}} * \prod_{Part \in \overline{part_{\mathcal{E}}}} |CT|_{Part}|, \quad (\text{A})$$

where  $\overline{part_{\mathcal{E}}} = \mathcal{P} \setminus part_{\mathcal{E}}$ . Note that when  $\overline{part_{\mathcal{E}}}$  is empty, we abusively consider that  $\prod_{\emptyset} = 1$ .

The presentation now focuses on calculating  $\beta_{\mathcal{E}}$ . Two cases arise depending on the size of  $\mathcal{E}$ .

<sup>1</sup>We choose to store passing partial tuples rather than eliminated tuples as they are usually less numerous. In an optimized version this decision could be made dynamically based on view cardinality estimation.

## Atomic explanations.

We start with explanations  $\mathcal{E}$  containing only one condition  $c$  (Algorithm 2 lines 3-5), which we call *atomic explanations*. To find the number of eliminated partial tuples  $\beta_{\mathcal{E}}$ , we firstly compute the set of passing partial tuples w.r.t.  $c$ , which we store in the view  $V_c$ :

$$V_c = \begin{cases} \pi_{\{R_{id} | R \in Part\}}(\sigma_c[CT]_{Part}) & \text{if } part_{\mathcal{E}} = \{Part\} \\ \pi_{\{R_{id} | R \in Part_1 \cup Part_2\}}([CT]_{Part_1} \bowtie_c [CT]_{Part_2}) & \text{if } part_{\mathcal{E}} = \{Part_1, Part_2\} \end{cases}$$

Then,

$$\beta_{\mathcal{E}} = \prod_{Part \in part_{\mathcal{E}}} |CT|_{Part}| - |V_c| \quad (\text{B})$$

**EXAMPLE 4.3.** For  $c_2$ , we have  $part_{c_2} = \{Part_1, Part_2\}$ , so  $V_{c_2} = \pi_{R_{Id}, S_{Id}, T_{Id}}([CT]_{Part_1} \bowtie_{R.B=T.B} [CT]_{Part_2})$ . This results in  $|V_{c_2}| = 4$ , and by Eq. (B) we obtain  $\beta_{c_2} = |CT|_{Part_1}| * |CT|_{Part_2}| - |V_{c_2}| = 3 * 4 - 4 = 8$ . Since all partitions of  $\mathcal{P}$  are in  $part_{c_2}$ , applying Equ. (A) results in  $\alpha_{c_2} = \beta_{c_2} = 8$ . For  $c_3$ ,  $\beta_{c_3} = |CT|_{Part_2}| - |V_{c_3}| = 4 - 2 = 2$ , so  $\alpha_{c_3} = 3 * 2 = 6$ . Fig. 4 (second level) displays the process for all atomic explanations.

## Non atomic explanations.

Now, consider  $\mathcal{E} = \{c_1, \dots, c_n\}$ ,  $n > 1$  (Alg. 2, lines 6-16). For the moment, we assume that the conditions in  $\mathcal{E}$  share the same schema, so the intersection and union of  $V_{c_i}$  for  $i = 1, \dots, n$  are well-defined. Firstly, we compute the view  $V_{\mathcal{E}}$  storing the passing partial tuples wrt  $\mathcal{E}$  as  $V_{\mathcal{E}} = V_{c_1} \cap \dots \cap V_{c_n}$ . To compute the number of partial tuples pruned out by  $\mathcal{E}$ , we need to find the number of partial tuples pruned out by  $c_1$  and  $\dots$  and  $c_n$ , i.e.,  $\beta_{\mathcal{E}} = |\overline{V_{c_1}} \cap \dots \cap \overline{V_{c_n}}|$ . By the well-known *DeMorgan law* [29], we have  $\beta_{\mathcal{E}} = |\overline{V_{c_1} \cup \dots \cup V_{c_n}}|$ , which spares us from computing the complements of  $V_{c_i}$ .

To compute the cardinality of the union among  $V_{c_i}$ , we rely on the *Principle of Inclusion and Exclusion for counting* [16]:

$$|\bigcup_{i=1}^n V_{c_i}| = \sum_{\emptyset \neq J \subseteq [n]} (-1)^{|J|+1} \left| \bigcap_{j \in J} V_{c_j} \right|$$

We further rewrite the previous formula to re-use results obtained for sub-combinations of  $\mathcal{E}$ , obtaining Eq. (C).

$$\begin{aligned} |\bigcup_{i=1}^n V_{c_i}| &= |\bigcup_{i=1}^{n-1} V_{c_i}| + |V_{c_n}| \\ &\quad + \sum_{\emptyset \neq J \subseteq [n-1]} (-1)^{|J|} \left| \bigcap_{j \in J} V_{c_j} \cap V_{c_n} \right| \end{aligned} \quad (\text{C})$$

At this point, we have all the necessary data to compute  $\beta_{\mathcal{E}}$ . However, so far we assumed that the conditions in  $\mathcal{E}$  have the same schema. In the general case, we have to “extend” the schema of a view  $V_c$  to the one of  $V_{\mathcal{E}}$ , in order to have well-defined set operations. The cardinality of an extended  $V_c^{ext}$  is given by Eq. (D).

$$|V_c^{ext}| = \prod_{Part \in part_{\mathcal{E}} \setminus part_c} |CT|_{Part}| * |V_c| \quad (\text{D})$$

Based on Eq. (D) we obtain Eq. (E) that generalizes Eq. (C).

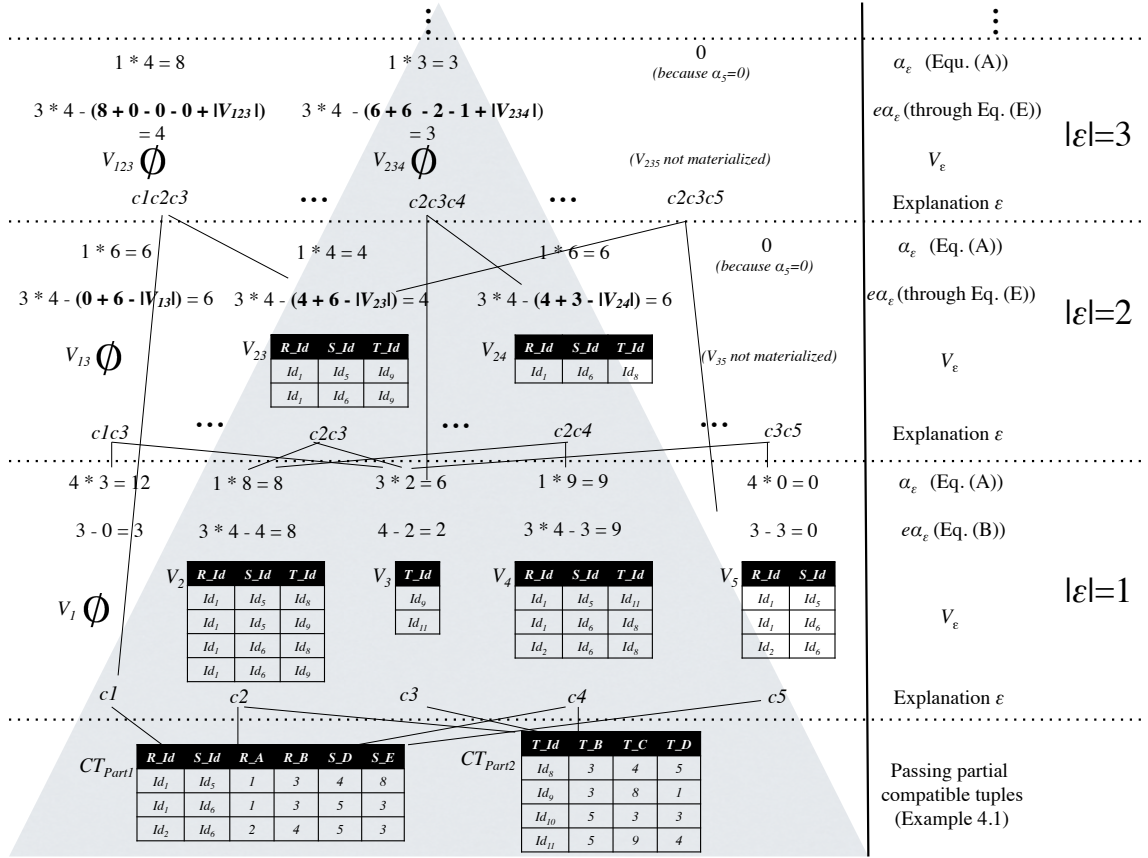


Figure 4: Running example with the different steps of *Ted++* (up to explanations of size 3) in Alg. 1 and Alg. 2

$$\begin{aligned}
 \left| \left( \bigcup_{i=1}^n V_{ci} \right)^{ext} \right| &= \left| \left( \bigcup_{i=1}^{n-1} V_{ci} \right)^{ext} \right| + |V_{cn}^{ext}| \\
 &+ \sum_{\emptyset \neq J \subseteq [n-1]} (-1)^{|J|} \left| \left( \bigotimes_{j \in J} V_{cj} \right) \bowtie V_{cn} \right|^{ext}
 \end{aligned} \quad (E)$$

In Eq. (E) we have replaced the intersection with natural join. The cardinalities of the views  $V_{\mathcal{E}'} = (\bigotimes_{j \in J} V_{cj})_{j \in J} \bowtie V_{cn}$  associated with  $\mathcal{E}'$  for  $|J| < n-1$ , have already been computed by previous steps and have only to be extended to the schema of  $V_E$ . When  $|J|=n-1$ , then  $V_{\mathcal{E}'}=V_E$ . A detailed discussion on how and when we materialize the view  $V_E$  is given shortly after.

Using the above, we trivially compute the number  $\beta_E$  of eliminated partial tuples as the complement of  $\left| \left( \bigcup_{i=1}^n V_{ci} \right)^{ext} \right|$  (Alg. 2, line 17). The number of eliminated tuples is then calculated by Eq. (A).

**EXAMPLE 4.4.** To illustrate the concepts introduced above, please follow on Fig. 4 below discussion.

For the explanation  $c_2c_3$ , Eq. (E) gives:  $|(V_2 \cup V_3)^{ext}| = |V_2^{ext}| + |V_3^{ext}| - |(V_2 \bowtie V_3)^{ext}|$ . The schema of  $part_{23} = \{Part_1, Part_2\}$  is  $\Gamma_{23} = \{R\_Id, S\_Id, T\_Id\}$ . The view  $V_2$  has already a matching schema, thus  $|V_2^{ext}| = |V_2| = 4$ . For  $V_3$ ,  $\Gamma_3 = \{T\_Id\}$ , we thus apply Eq. (D) and obtain  $|V_3^{ext}| = |V_{Part1}| * |V_3| = 3 * 2 = 6$ . Still,  $|V_{23}| = |(V_2 \bowtie V_3)^{ext}|$  remains to be calculated. Intuitively, because  $V_2$  and  $V_3$  target schemas share attribute  $T\_Id$ ,  $V_{23} = V_2 \bowtie_{T\_Id} V_3$ . The view  $V_{23}$  is materialized and contains 2 tuples (as shown in Fig. 4). So, finally, from Eq. (E) we obtain

$|(V_2 \cup V_3)^{ext}| = 4 + 6 - 2 = 8$ . Since  $|Part_1| * |Part_2| = 12$  then  $\beta_{23} = 12 - 8 = 4$ , and by Eq. (A)  $\alpha_{23} = 4$ .

We now focus on the explanation  $c_3c_5$ . The schemas of  $V_3$  and  $V_5$  are disjoint and intuitively  $V_{35} = V_3 \times V_5$ . Here,  $V_{35}$  is not materialized, we simply calculate  $|V_{35}| = |V_3| * |V_5| = 6$ . Then,  $|\beta_{35}| = 12 - (12 + 6 - 6) = 0$ . As we will see later, these steps are never performed in our algorithm. The fact that  $c_5$  does not eliminate any tuple (see  $\alpha_5 = 0$  in Fig. 4) implies that neither do any of its super-combinations. Thus, a priori we know that  $\alpha_{35} = \alpha_{235} = \dots = 0$ .

Finally, we illustrate the case of a bigger size combination, for example  $c_2c_3c_4$  of size 3. Eq. (E) yields  $|(V_2 \cup V_3 \cup V_4)^{ext}| = |(V_2 \cup V_4)^{ext}| + |V_3^{ext}| - |(V_2 \bowtie V_3)^{ext}| - |(V_4 \bowtie V_3)^{ext}| + |(V_2 \bowtie V_3 \bowtie V_4)^{ext}|$ . All terms of the right side of the equation are available from previous iterations, except for  $|(V_2 \bowtie V_3 \bowtie V_4)^{ext}|$ . As before, we check the common attributes of the views and obtain  $V_{234} = V_{24} \bowtie_{R\_Id, S\_Id, T\_Id} V_3$ . So,  $|(V_2 \cup V_3 \cup V_4)^{ext}| = 6 + 6 - 2 - 1 + 0 = 9$  and  $\beta_{234} = \alpha_{234} = 12 - 9 = 3$ . In the same way, we compute all the possible explanations until  $c_1c_2c_3c_4c_5$ .

#### View Materialization: when and how.

To decide when and how to materialize the views for the explanations, we partition the set of the views associated with the conditions in  $\mathcal{E}$ . Consider the relation  $\sim$  defined over these views by  $V_i \sim V_j$  if the target schemas of  $V_i$  and  $V_j$  have at least one common attribute. Consider the transitive closure  $\sim^*$  of  $\sim$  and the induced partitioning of  $V_E$  through  $\sim^*$ .

When this partitioning is a singleton,  $V_E$  needs to be materialized (Alg. 2, line 9). The materialization of  $V_E$  is specified by joining the views associated with the sub-conditions, which may



be done in more than one way, as usual. For example, for the combination  $c_2c_3c_4$ ,  $V_{234}$  can either be computed through  $V_{23} \bowtie V_4$  or  $V_{24} \bowtie V_3$  or  $V_{34} \bowtie V_2$  or  $V_2 \bowtie V_3 \bowtie V_4$ ... because all these views are known from previous iterations. The choice of the query used to materialize  $V_E$  is done based on a cost function. This function gives priority to materializing  $V_E$  by means of one join, which is always possible: because  $V_E$  needs to be materialized, we know that at least one view associated with a sub-combination of size  $n-1$  has been materialized. In other words, priority is given to using at least one materialized view associated with one of the largest sub-combinations. For our example, it means that either  $V_{23} \bowtie V_4$  or  $V_{24} \bowtie V_3$  or  $V_{34} \bowtie V_2$  is considered. In order to choose among the one-join queries computing  $V_E$ , we favor a one-join query  $V_i \bowtie V_j$  minimal w.r.t.  $|V_i| + |V_j|$ . For the example, and considering also Fig. 4 we find that  $|V_2| + |V_{34}| = |V_4| + |V_{23}| = 5$  and  $|V_3| + |V_{24}| = 3$ . So, the query used for the materialization is  $V_3 \bowtie V_{24}$  (its result being empty in our example). Nevertheless, we avoid the materialization of  $V_E$  if the partitioning is a singleton (Alg. 2, line 9 & 16), when for some sub-combination  $\mathcal{E}'$  of  $\mathcal{E}$  it was computed that  $\alpha_{\mathcal{E}'} = 0$ . In that case, we know a priori that  $\alpha_{\mathcal{E}} = 0$  (see Ex. 4.4).

If the partitioning is not a singleton,  $V_E$  is not materialized (Alg. 2, line 14). For example, the partitioning for  $c_3c_5$  is not a singleton and so the size  $|V_{35}| = |V_3| \times |V_5| = 6$ .

### Post-processing.

In Alg. 2 we associated with each possible explanation  $\mathcal{E}$  the number of eliminated tuples  $\alpha_{\mathcal{E}}$ . However, recall that the calculation of this number so far counts any tuple eliminated by  $\mathcal{E}$ , even though the same tuples may be eliminated by some super-combinations of  $\mathcal{E}$  (see Ex. 4.5). This means that for some tuples, multiple explanations have been assigned. To make things even, the last step of *Ted++* (Alg. 1, line 6) is about calculating the coefficient of  $\mathcal{E}$  by subtracting the coefficients of its super-combinations from  $\alpha_{\mathcal{E}}$ :

$$coef_{\mathcal{E}} = \alpha_{\mathcal{E}} - \left( \sum_{\mathcal{E}' \subseteq \mathcal{E}} coef_{\mathcal{E}'} \right) \quad (F)$$

EXAMPLE 4.5. Consider known  $coef_{1234} = 2$  and  $coef_{123} = 2$ . We have found in Ex. 4.4 that  $\alpha_{23} = 4$ . With Eq. (F),  $coef_{23} = 4 - 2 - 2 = 0$ . In the same way  $coef_2 = 4 - 0 - 2 - 2 = 0$ . The algorithm leads to the expected Why-Not explanation polynomial already provided in Ex. 3.4.

### 4.3 Complexity analysis.

In the pseudo-code for *Ted++* provided in Alg. 1, we can see that *Ted++* divides into the phases of (i) partitioning  $\mathcal{S}$ , (ii) materializing a view for each partition, (iii) computing the explanations, and (iv) computing the exact coefficients. When computing the explanations, according to Alg. 2, *Ted++* iterates through  $2^{|C|}$  condition combinations and for each, it decides upon view materialization (again through partitioning) before materializing it, or simply calculates  $|V_E|$  before applying equations to compute  $\alpha_{\mathcal{E}}$ . Overall, we consider that all mathematical computations are negligible so, the worst case complexities of steps (i) through (iv) are  $O(|\mathcal{S}| + |WN|) + O(|\mathcal{S}|) + O(2^{|C|}(|\mathcal{S}| + |C|)) + O(2^{|C|})$ . For large enough queries, we can assume that  $|\mathcal{S}| + |C| < 2^{|C|}$ , in which case the complexity simplifies to  $O(2^{|C|})$ .

Obviously, the complexity analysis above does not take into account the cost of actually materializing views; in its simplified form, it only considers how many views need to be materialized in the worst case. Assume that  $n = \max(|I_R| |R \in \mathcal{S}|)$ . The materialization of any view is bound by the cost of materializing a cross product over the relations involved in the view - in the worst

Table 2: Queries for the scenarios in Tab. 3

Query	Expression
Q1	$C \bowtie_{sector} W \bowtie_{witnessName} S \bowtie_{hair,clothes} P$
Q2	$\sigma_{C.sector > 99}[C] \bowtie_{sector} W \bowtie_{witnessName} S \bowtie_{hair,clothes} P$
Q3	$W \bowtie_{sector2} C2 \bowtie_{sector1} \sigma_{C.type=Aiding}[C]$
Q4	$P2 \bowtie_{name,hair} \sigma_{P1.name < B}[P1]$
Q5	$L \bowtie_{movieId} \sigma_{M.year > 2009}[M] \bowtie_{name} \sigma_{R.rating > 8}[R]$
Q6	$\sigma_{AA.party=Republican}[AA] \bowtie_{id} \sigma_{Co.Byear > 1970}[Co]$
Q7	$E \bowtie_{eId} \sigma_{ES.sub=Sen.Com.}[ES] \bowtie_{id} \sigma_{SPO.party=Rep.}[SPO]$
Qs3	$\sigma_{type=Aiding}[Q2]$
Qs4	$\sigma_{witnessname > S}[Qs3]$
Qj	$C \bowtie_{sector} \sigma_{name > S}[W]$
Qj2	$Q_j \bowtie_{witnessname} S$
Qj3	$Q_{j2} \bowtie_{clothes} P$
Qj4	$Q_{j3} \bowtie_{hair} P$
Qc	$L \bowtie_{id} L2 \bowtie_{M2.mid=L2.mid} M2 \bowtie_{year,mid} \sigma_{year=1980}[M1]$
Qtpch	$C \bowtie_{ckey} \sigma_{odate < 1998-07-21}[O] \bowtie_{okey} \sigma_{sdate > 1998-07-21}[L]$

case  $O(n^{|\mathcal{S}|})$ . This yields a combined complexity of  $O(2^{|C|} n^{|\mathcal{S}|})$ . However, *Ted++* in the general case (more than one induced partitions), has a tighter upper bound:  $O(n^{k_{\mathcal{E}1}} + n^{k_{\mathcal{E}2}} + \dots + n^{k_{\mathcal{E}N}})$ , where  $k_{\mathcal{E}} = |\text{part}_{\mathcal{E}}|$ , for all combinations  $\mathcal{E}$  and  $N = 2^{|C|}$ . It is easy to see that  $n^{k_{\mathcal{E}1}} + n^{k_{\mathcal{E}2}} + \dots + n^{k_{\mathcal{E}N}} < 2^{|C|} n^{|\mathcal{S}|}$ , when there is more than one partition.

## 5. EXPERIMENTAL EVALUATION

This section presents an experimental evaluation of *Ted++*, using real and synthetic datasets. In Sec. 5.1, we compare *Ted++* with the existing algorithms returning query-based explanations, i.e., with NedExplain [5] and Why-Not [9]. Sec. 5.2 studies the runtime of *Ted++* with respect to various parameters that we vary in a controlled manner. We have implemented *Ted*, *Ted++*, NedExplain, and Why-Not in Java. We ran the experiments on a Mac Book Air, running MAC OS X 10.9.5 with 1.8 GHz Intel Core i5, 4GB memory, and 120GB SSD. We used PostgreSQL 9.3 as database system.

### 5.1 Comparative Evaluation

The comparative evaluation to Why-Not and NedExplain considers both efficiency (runtime) and effectiveness (explanation quality). When considering efficiency, we also include *Ted* in the comparison (*Ted* producing the same Why-Not explanation as *Ted++*).

For the experiments in this section, we have used data from three databases named *crime*, *imdb*, and *gov*. The *crime* database corresponds to the sample crime database of the Trio system (available at <http://infolab.stanford.edu/trio/>) and was previously used to evaluate Why-Not and NedExplain. The data describes crimes and involved persons (suspects and witnesses). The *imdb* database contains real-world movie data from IMDB (<http://www.imdb.com>). Finally, the *gov* database contains information about US congressmen and financial activities (data from <http://bioguide.congress.gov>, <http://usaspending.gov>, and <http://earmarks.omb.gov>).

For each dataset, we have created a series of scenarios (crime1-gov5 in Tab. 3 - ignore remaining scenarios for now). Each scenario consists of a query further defined in Tab. 2 (Q1-Q7) and a simple Why-Not question, as Why-Not and NedExplain support only this type of Why-Not question. The queries have been designed to include queries with a small set of conditions (Q6) or a larger one (Q1, Q3, Q5, Q7), containing self-joins (Q3, Q4), having empty intermediate results (Q2), as well as containing inequalities (Q2, Q4, Q5, Q6).

#### 5.1.1 Why-Not Explanation Evaluation

In Tab. 1 we report that the explanations returned by Why-Not

**Table 3: Scenarios**

Scenario	Query	Why-Not question
crime1	Q1	{P.Name=Hank,C.Type=Car theft}
crime2	Q1	{P.Name=Roger,C.Type=Car theft}
crime3	Q2	{P.Name=Roger,C.Type=Car theft}
crime4	Q2	{P.Name=Hank,C.Type=Car theft}
crime5	Q2	{P.Name=Hank}
crime6	Q3	{C2.Type=kidnapping}
crime7	Q3	{W.Name=Susan,C2.Type=kidnapping}
crime8	Q4	{P2.Name=Audrey}
imdb1	Q5	{name=Avatar}
imdb2	Q5	{name=Christmas Story,L.locationId=USANew York}
gov1	Q6	{Co.firstname=Christopher}
gov2	Q6	{Co.firstname=Christopher,Co.lastname=MURPHY}
gov3	Q6	{Co.firstname=Christopher,Co.lastname=GIBSON}
gov4	Q7	{sponsorId=467}
gov5	Q7	{SPO.sponsorIn=Lugar,E.camout>=1000}
crime <sub>s</sub> - crime <sub>s4</sub>	Q1,Q2, Q <sub>s3</sub> ,Q <sub>s4</sub>	{P.Name=Hank,C.Type=Car theft}
crime <sub>j</sub> - crime <sub>j4</sub>	Q <sub>j</sub> - Q <sub>j4</sub>	{W.name=Jane, C.type=Car theft}
imdb <sub>c</sub>	Q <sub>c4</sub>	{L2.locationId=L1.locationId, M1.mid=L2.mid, L1.year>L2.year,M1.name=Duck Soup}
imdb <sub>c2</sub>	Q <sub>c4</sub>	{(L2.locationId=L1.locationId, M1.mid=L2.mid, L1.year>L2.year)}
crime <sub>5c2</sub>	Q2	{P.Name=Hank, C.type=Car theft}
crime <sub>5c3</sub>	Q2	{P.Name=Hank, C.type=Car theft, S.witness=Aphrodite}
crime <sub>5c4</sub>	Q2	{(P.Name=Hank, C.type=Car theft, S.witness=Aphrodite, W.sector=34)}
crime <sub>5c5</sub>	Q2	{P.Name=Hank, C.type=Car theft, S.witness=Aphrodite, W.sector=34,S.hair=green}
imdb <sub>cc</sub>	Q <sub>c</sub>	{M.year>M2.year}
tpch <sub>s</sub>	Q <sub>tpch</sub>	{L.extprice>50000,O.odate<1996-01-01}
tpch <sub>c</sub>	Q <sub>tpch</sub>	{L.extprice>100000,O.odate=L.cdate, C.nkey=4}

and NedExplain consist of sets of query conditions, whereas *Ted++* returns a polynomial of query conditions. For comparison purposes, we trivially map *Ted++*'s Why-Not explanation to sets of conditions, e.g.,  $3c_3 * c_4 + 2c_3 * c_6$  maps to  $\{\{c_3, c_4\}, \{c_3, c_6\}\}$ . For conciseness, we abbreviate condition sets, e.g., to  $c_{34}, c_{36}$ .

Tab. 4 summarizes the Why-Not explanations of the three algorithms. These scenarios make apparent that the explanations by NedExplain or Why-Not are incomplete, in two senses. First, they produce only a subset of the possible explanations, failing to provide alternatives that could be useful to the user when she tries to fix the query. Second, even the explanation they provide may lack parts, which can drive the user to fruitless fixing attempts. On the contrary, *Ted++* produces all the possible, complete explanations.

For the first argument, consider the scenario *gov2*. Why-Not and NedExplain return  $c_1$  and  $c_3$  respectively, but they both fail to indicate that both the explanations are valid, as opposed to *Ted++*. Then, consider *crime8*. NedExplain returns the join  $c_2$  ( $S \bowtie_{hair} P$ ) - Why-Not falsely does not produce any explanations in this case. *Ted++* indicates that except for this join, the selection  $c_3$  ( $\sigma_{name < 'B'}[P]$ ) for instance is also an explanation. From a developer's perspective, selections are typically easier or more reasonable to change. So, having the complete set of explanations potentially provides the developer with useful alternatives.

For the second argument consider *crime5*. NedExplain returns  $c_1$  ( $C \bowtie_{sector} W$ ). The explanation of *Ted++* does not contain the atomic explanation  $c_1$ , but there exist combinations including  $c_1$  as a part, like  $c_{15}$ . This means that the explanation by NedExplain is incomplete; a repair attempt of  $c_1$  alone will never yield the desired results. Similarly, *crime7* illustrates a case, when the Why-Not algorithm produces an explanation ( $c_3$ ) that misses some parts. Then, in *gov3* NedExplain and Why-Not both return  $c_2$ . However, let us now assume the developer prefers to not change this condition. Keeping in mind that those algorithms' answers may change

**Table 4: *Ted++*, Why-Not, NedExplain answers per scenario**

Scenario	<i>Ted++</i>	Why-Not	NedExplain
crime1	$c_{1234}, \dots, c_{12}, c_3, c_2, c_1$		$c_1$
crime2	$c_{1234}, c_{34}, c_{13}, \dots, c_3$	$c_{34}$	$c_{34}, c_1$
crime3	$c_{12345}, \dots, c_{145}, c_{345}, c_{35}$	$c_{34}, c_5$	$c_5, c_{34}$
crime4	$c_{12345}, \dots, c_{25}, c_{15}$	$c_5$	$c_1, c_5$
crime5	$c_{12345}, \dots, c_{15}, c_5$	$c_5$	$c_1$
crime6	$c_{123}, c_{31}, c_{23}, c_{12}, c_3, c_2, c_1$	$c_3$	$c_2$
crime7	$c_{123}, c_{13}, c_{12}, c_1$	$c_3$	$c_2, c_1$
crime8	$c_{23}, c_3, c_2, c_1$		$c_2$
imdb1	$c_{123}, c_{13}, c_{23}, c_3$	$c_3$	$c_3, c_2$
imdb2	$c_{13}$		$c_1, c_3$
gov1	$c_{123}, c_{13}, c_{23}, c_{12}, c_3, c_2, c_1$	$c_3$	$c_2, c_3$
gov2	$c_{13}, c_3, c_1$	$c_1$	$c_3$
gov3	$c_{123}, c_{23}, c_2$	$c_2$	$c_2$
gov4	$c_{123}, c_{23}, c_2$	$c_3$	$c_3, c_2$
gov5	$c_{124}, c_{14}, c_{24}, c_{12}, c_4, c_2, c_1$	$c_1$	$c_1$

when changing the query tree, she may start trying different trees to possibly obtain a Why-Not explanation without  $c_2$ . Knowing the explanation of *Ted++* prevents her from spending any effort on this, as it shows that all explanations contain  $c_2$  as a part.

By mapping the explanation of *Ted++* to sets of explanations, we have let aside an important property: the coefficients of the polynomial. For example, the complete Why-Not explanation polynomial of *crime8* is  $2384 * c_{23} + 20 * c_3 + 4 * c_1 + 8 * c_2$ . Assume that the developer would like to recover at least five missing tuples, by changing as few conditions as possible. The polynomial suggests to change either  $c_3$  or  $c_2$ : they both require one condition change and provide the possibility of obtaining up to 20 and 8 missing tuples, respectively. The coefficient of  $c_1$  being 4, does not make  $c_1$  a good candidate, whereas  $c_2c_3$  require two condition changes. Clearly, the results of NedExplain or Why-Not are not informative enough for such a discussion.

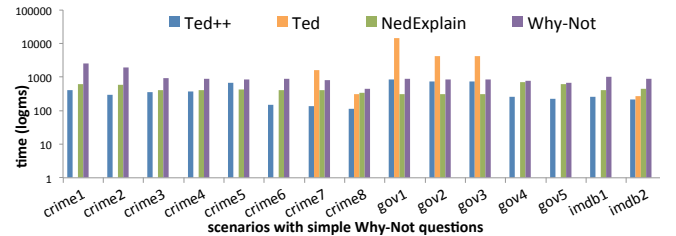
### 5.1.2 Runtime Evaluation

We now compare the performance w.r.t. runtime of *Ted++* with the other algorithms.

#### *Ted++* vs. NedExplain and Why-Not.

For this comparative evaluation, we again consider scenarios *crime1* through *gov5* of Tab. 3 as they involve simple Why-not questions, making them processable by all three algorithms. Fig. 5 summarizes the runtimes in logarithmic scale for each algorithm and scenario. We observe that the runtime of *Ted++* is always comparable to the runtime of NedExplain and that in some cases, it is significantly faster than Why-Not.

Why-Not traces compatible tuples based on tuple lineage stored in Trio. As already stated in [5, 9], this design choice slows down Why-Not performance. On the contrary, both NedExplain and *Ted++* compute the compatible data more efficiently by issu-


**Figure 5: Runtimes for *Ted++*, *Ted*, NedExplain and Why-Not**

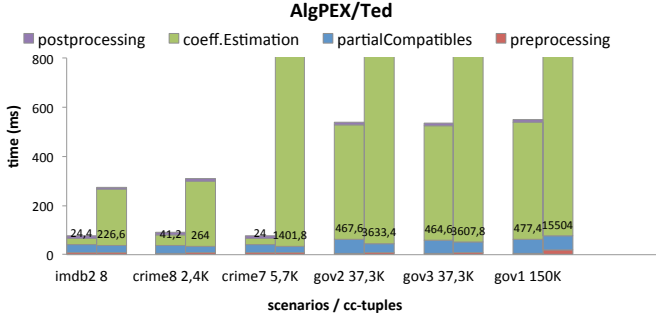


Figure 6: *Ted++* and *Ted* runtime distribution

ing SQL statements to the database and further using the unique identifiers of the source tuples. We claim that a better implementation choice for tuple tracing in Why-Not would yield a runtime comparable to NedExplain, a claim backed up by their comparable runtime complexities. Another definition and implementation issue of both Why-Not and NedExplain, which explains the sometimes faster runtime of *Ted++* is the fact that their input is potentially much larger as it includes the full database instance instead of the compatible data only. Clearly, this slows the tracing of compatible data through the query tree.

Let us see what happens when *Ted++* is slower than - but still comparable to - NedExplain, for example in *gov1-gov3*. In NedExplain all compatible tuples are pruned out by conditions very close to the leaf level of the query tree, so the bottom-up traversal of the tree can stop very early. *Ted++* always “checks” all conditions so cannot benefit from such an early termination. However, this runtime improvement of NedExplain often comes at the price of incomplete explanations (e.g., *gov1*).

*Ted++* vs. *Ted*.

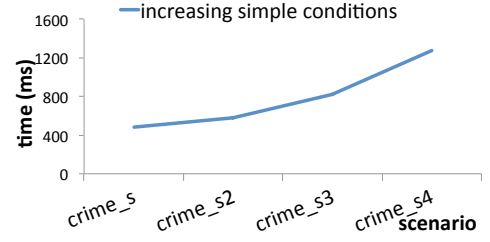
Fig. 5 also reports runtimes for *Ted* on 6 out of 15 scenarios (for the others, *Ted* runs out of time). To experimentally demonstrate where *Ted*’s problem lies, we compare the time distribution in *Ted* and *Ted++*.

Fig. 6 divides the runtime into four common phases. Among these, the *coefficientEstimation* phase is the one that is inherently different in both algorithms. *Ted* iterates over the set of compatible tuples and computes the explanation for each one. *Ted++* explores the search space of possible explanations and calculates, based on the number of passing partial compatible tuples, the number of compatible tuples eliminated by each explanation. Thus, this is the phase in which we expect to have an important runtime difference between *Ted* and *Ted++*. In reporting the phase-wise runtime, Fig. 6 cuts the bar for *Ted* in the scenarios *crime7*, *gov1*, *gov2* and *gov3* as the execution time is much higher compared to the other scenarios and to the runtime of *Ted++* (the runtime of the *coefficientEstimation* phase is the label on the respective bars).

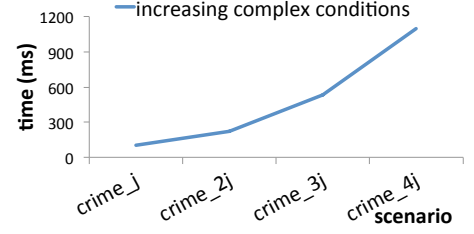
As said in Sec. 4, *Ted*’s main issue w.r.t. efficiency is its strong dependence on the number of compatible tuples. This is experimentally observed in Fig. 6: with the growth of the set of compatible tuples in the scenarios, the time dedicated to *coefficientEstimation* also grows (the scenarios are reported in an ascending order). *Ted++* depends on the number of compatible tuples as well, but not as strongly as *Ted*. This can be seen in *crime8* and *crime7*, or *gov3* and *gov1*; while the number of tuples grows, *Ted++*’s *coefficientEstimation* phase remains roughly steady.

## 5.2 *Ted++* Analysis

We now study *Ted++*’s behavior w.r.t. the following parameters:



(a) simple conditions



(b) complex conditions

Figure 7: *Ted++* runtime w.r.t. number of conditions in  $q$

(i) the type (simple or complex) of the input query  $Q$  and the number of  $Q$ ’s conditions, (ii) the type of the Why-Not question (simple or complex) and the number and selectivity of conditions the Why-Not question involves, and (iii) the size of the database instance  $\mathcal{I}$ . Note that (ii) and (iii) are tightly connected with the number of compatible tuples, which is one of the main parameters influencing the performance. In addition to the number of compatible tuples, another important factor is the selectivity of the query conditions over the compatible data (i).

### Experimental Setup.

For the parameter variations (i) and (ii), we use again the *crime*, *imdb*, and *gov* databases. To adjust the database instance size for case (iii), we use data produced by the TPC-H benchmark data generator (<http://www.tpc.org/tpch/>). More specifically, we generate instances of 1GB and 10GB and further produce smaller data sets of 10MB and 100MB to obtain a series of datasets whose size differs by a factor of 10. In this paper, we report results for the original query Q3 of the TPC-H set of queries. It includes two complex and three simple conditions, two of which are inequality conditions. Since the original TPC-H query Q3 is an aggregation query, we have changed the projection condition. The queries used in this section are summarized in Tab. 2 ( $Q_s$ - $Q_{tpch}$ ) and the scenarios in Tab. 3 (*crime<sub>s</sub>-tpch<sub>c</sub>*).

### Adjusting the query.

Given a fixed database instance and Why-Not question, we start from query Q1 and gradually add simple conditions, yielding the series of queries Q1, Q2,  $Q_{s3}$ ,  $Q_{s4}$ . The evolution of *Ted++* runtime for this series of queries is shown in Fig. 7 (a). Similarly, starting from query  $Q_j$ , we introduce step by step complex conditions, yielding  $Q_j$ - $Q_{j4}$ . Corresponding runtime results are reported in Fig. 7 (b).

As expected, in both cases, increasing the number of query conditions (either complex or simple) results in increasing runtime. The incline of the curve depends on the selectivity of the introduced condition; the less selective the condition the steeper the line becomes. This is easy to explain, as in the *coefficientEstimation* phase, a view contains more tuples (=passing partial tuples) when the condition is less selective. This results in more computations in the super-combinations iterations, leaving space for further opti-

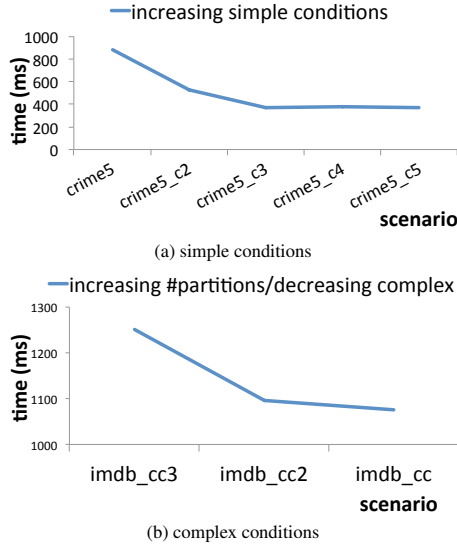


Figure 8: *Ted++* runtime w.r.t. number of conditions in *WN*

mization by dynamically deciding on passing vs eliminated tuples materialization.

Note that the curve in Fig. 7 (a) starts at a much higher point than in Fig. 7 (b). This is because the query  $Q_1$  (*crime<sub>s</sub>*) initially includes four complex conditions, in contrast to  $Q_j$  (*crime<sub>10</sub>*) that includes one complex and one simple condition.

#### Adjusting the Why-Not question.

Next, we vary the type and the number of conditions in the Why-Not question *WN*. Fig. 8 shows the cases when we start (a) with a simple *WN* and progressively add more simple conditions and (b) start with a complex *WN* and progressively add more complex conditions.

The scenarios considered for Fig. 8 (a) have as starting point the simple scenario *crime<sub>5</sub>* (see Tab. 3). Then, keeping the same input instance and query, we add attribute-constant comparisons to *WN*, a procedure resulting in fewer tuples in each step. As expected, the more conditions (the less tuples) the faster the Why-Not explanation is returned, until we reach a certain point (here from *crime<sub>5</sub><sub>c3</sub>* on). From this point, the runtime is dominated by the time to communicate with the database that is constant over all scenarios.

As we introduce complex conditions in the *WN*, the number of generated partitions (potentially) drops as more relations are included in a same partition. To study the impact of the induced number of partitions in isolation, we keep the number of the compatible tuples constant in our series of complex scenarios (*imdb<sub>cc</sub>*-*imdb<sub>cc3</sub>*). The number of partitions entailed by *imdb<sub>cc</sub>*, *imdb<sub>cc2</sub>*, and *imdb<sub>cc3</sub>* are 3, 2, and 1, respectively. The results of Fig. 8 (b) confirm our theoretical complexity discussion, i.e., as the number of partitions decreases, the time needed to produce the Why-Not explanation increases.

#### Increasing size of input instance.

The last parameter we study is the input database size. To this end, we have created two scenarios, one with a simple and one with a complex Why-Not question *WN*, and both using the same query  $Q_{tpch}$ . We run both scenarios for database sizes 10MB, 100MB, 1GB, and 10GB. The simple *WN* includes two inequality conditions, in order to be able to compute a reasonable number of compatible tuples. The complex *WN* contains one complex condition,

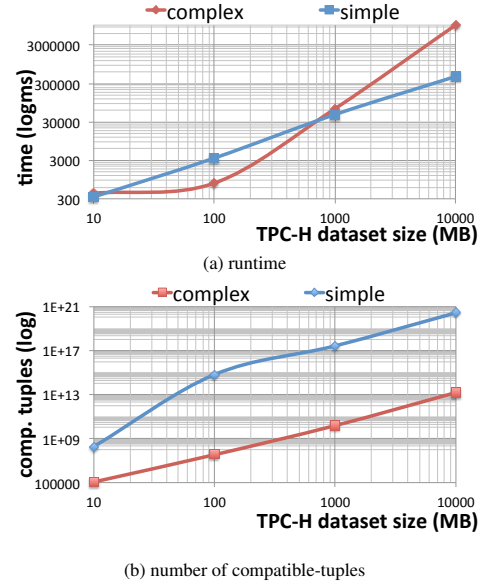


Figure 9: *Ted++* (a) runtime, and (b) number of compatible tuples for increasing database size, complex and simple *WN*

one inequality simple condition and one equality simple condition. It thus represents an average complex Why-Not question, creating two partitions over three relations.

Fig. 9 (a) shows the runtimes for both scenarios. The increasing runtime is tightly coupled to the fact that the number of computed tuples is augmenting proportionally to the database size, as shown in Fig. 9 (b). We observe that for small datasets (<500MB) in the complex scenario *Ted++*'s performance decreases with a low rate, whereas the rate is higher for larger datasets. For the simple scenario, runtime deteriorates in a steady pace. This behavior is aligned with the theoretical study; when the number of partitions is decreasing the complexity rises.

In summary, our experiments have shown that *Ted++* generates a more informative, useful and complete Why-Not explanation than the state of the art. Moreover, *Ted++* is competitive in terms of runtime. The dedicated experimental evaluation on *Ted++* verifies that it can be used in a large variety of scenarios with different parameters and that the obtained runtimes match the theoretical expectations. Finally, the fact that the experiments were conducted on an ordinary laptop, with no special capabilities in memory or disk space, supports *Ted++*'s feasibility.

## 6. CONCLUSION AND OUTLOOK

This paper provides a framework for Why-Not explanations based on polynomials, which enables to consider relational databases under set, bag and probabilistic semantics in a unified way. To efficiently compute the Why-Not explanation polynomial under set semantics we have designed a new algorithm *Ted++*, whose main feature is to completely avoid enumerating and iterating over the set of compatible tuples, thus it significantly reduces both space and time consumption. Our experimental evaluation showed that *Ted++* is at least as efficient as existing algorithms while providing useful insights in its Why-Not explanation for a developer. Also, we saw that *Ted++* scales well with various parameters, making it a practical solution. The proposed Why-Not explanation polynomial are easy to extend for unions of conjunctive queries, whereas an extension is not trivial for aggregation queries and is subject to future work.

Currently, we have been working on exploiting the Why-Not explanation polynomial to efficiently rewrite a query in order to include the missing answers in its result set. As there are many rewriting possibilities, we plan to select the most promising ones based on a cost function, built with the polynomial. For instance, we may rank higher rewritings with minimum condition changes (i.e., small combinations), minimum side-effects (i.e., small coefficients), etc.

## 7. REFERENCES

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [2] A. Baid, W. Wu, C. Sun, A. Doan, and J. F. Naughton. On debugging non-answers in keyword search systems. In *EDBT*, 2015.
- [3] N. Bidoit, M. Herschel, and K. Tzompanaki. Immutably answering why-not questions for equivalent conjunctive queries. In *TAPP*, 2014.
- [4] N. Bidoit, M. Herschel, and K. Tzompanaki. Immutably answering why-not questions for equivalent conjunctive queries. In *BDA*, 2014.
- [5] N. Bidoit, M. Herschel, and K. Tzompanaki. Query-based why-not provenance with NedExplain. In *EDBT*, 2014.
- [6] N. Bidoit, M. Herschel, and K. Tzompanaki. Efficient computation of polynomial explanations of why-not questions. In *CIKM (to appear)*, 2015.
- [7] D. Calvanese, M. Ortiz, M. Simkus, and G. Stefanoni. Reasoning about explanations for negative query answers in DL-Lite. *JAIR*, 2013.
- [8] B. t. Cate, C. Civili, E. Sherkhonov, and W.-C. Tan. High-level why-not explanations using ontologies. In *PODS*, 2015.
- [9] A. Chapman and H. V. Jagadish. Why not? In *SIGMOD*, 2009.
- [10] L. Chen, X. Lin, C. S. Jensen, and J. Xu. Answering why-not questions on spatial keyword top-k queries. In *ICDE*, 2015.
- [11] J. Cheney, L. Chiticariu, and W. C. Tan. Provenance in databases: Why, how, and where. *ACM-FTD*, 1(4), 2009.
- [12] Y. Cui, J. Widom, and J. L. Wiener. Tracing the lineage of view data in a warehousing environment. *TODS*, 25(2), 2000.
- [13] Y. Gao, Q. Liu, G. Chen, B. Zheng, and L. Zhou. Answering why-not questions on reverse top-k queries. *PVLDB*, 8(7):738–749, 2015.
- [14] T. J. Green, G. Karvounarakis, and V. Tannen. Provenance semirings. In *PODS*, 2007.
- [15] T. Grust and J. Rittinger. Observing sql queries in their natural habitat. *TODS*, 38(1):3, 2013.
- [16] M. Hall. *Combinatorial theory*, volume 71. John Wiley & Sons, 1998.
- [17] Z. He and E. Lo. Answering why-not questions on top-k queries. In *ICDE*, 2012.
- [18] M. Herschel. Wondering why data are missing from query results? ask Conseil Why-Not. In *CIKM*, 2013.
- [19] M. Herschel. A hybrid approach to answering why-not questions on relational query results. *ACM-JDIQ*, 5(3):10:1–10:29, 2015.
- [20] M. Herschel and M. A. Hernández. Explaining missing answers to SPJUA queries. *PVLDB*, 3(1), 2010.
- [21] J. Huang, T. Chen, A. Doan, and J. F. Naughton. On the provenance of non-answers to queries over extracted data. *PVLDB*, 1(1), 2008.
- [22] M. S. Islam, R. Zhou, and C. Liu. On answering why-not questions in reverse skyline queries. In *ICDE*, 2013.
- [23] N. Khousainova, Y. Kwon, M. Balazinska, and D. Suciu. SnipSuggest: Context-aware autocompletion for SQL. *PVLDB*, 4(1), 2010.
- [24] A. Meliou, W. Gatterbauer, K. F. Moore, and D. Suciu. The complexity of causality and responsibility for query answers and non-answers. *PVLDB*, 2011.
- [25] D. Mottin, A. Marascu, S. B. Roy, G. Das, T. Palpanas, and Y. Velegrakis. A probabilistic optimization framework for the empty-answer problem. *PVLDB*, 6(14), 2013.
- [26] A. Nandi and H. V. Jagadish. Guided interaction: Rethinking the query-result paradigm. *PVLDB*, 4(12), 2011.
- [27] S. Roy and D. Suciu. A formal approach to finding explanations for database queries. In *SIGMOD*, 2014.
- [28] Q. T. Tran and C.-Y. Chan. How to ConQueR why-not questions. In *SIGMOD*, 2010.
- [29] R. Vaught. *Set Theory An Introductory*. Birkhaeuser, 2001.